

## Design document

The implementation is divided into the following steps/jobs.

1. Count the number of nodes in the graph.
2. Get the number of edges in the graph, calculate minimum, maximum and average number of edges per node.
3. Prepare initial file with the format < NODE ID: updated rank: last rank :<adjacency list>>  
Initial rank is given as  $(1/\text{total nodes})$ , in the beginning both updated and last rank are set to this value.
4. Specify maximum number of iterations and slackness which defines the accuracy of convergence.
5. Calculate the contribution from dangling nodes.
6. Update the page-rank –  
Uses the file from previous iteration as input, takes into account dangling nodes.  
Makes use of the PageRank formula specified in paper with damping factor = .85.
7. Check if the ranks have converged.
8. Iterate steps 5,6,7 until convergence is reached or the maximum iteration limit is breached.

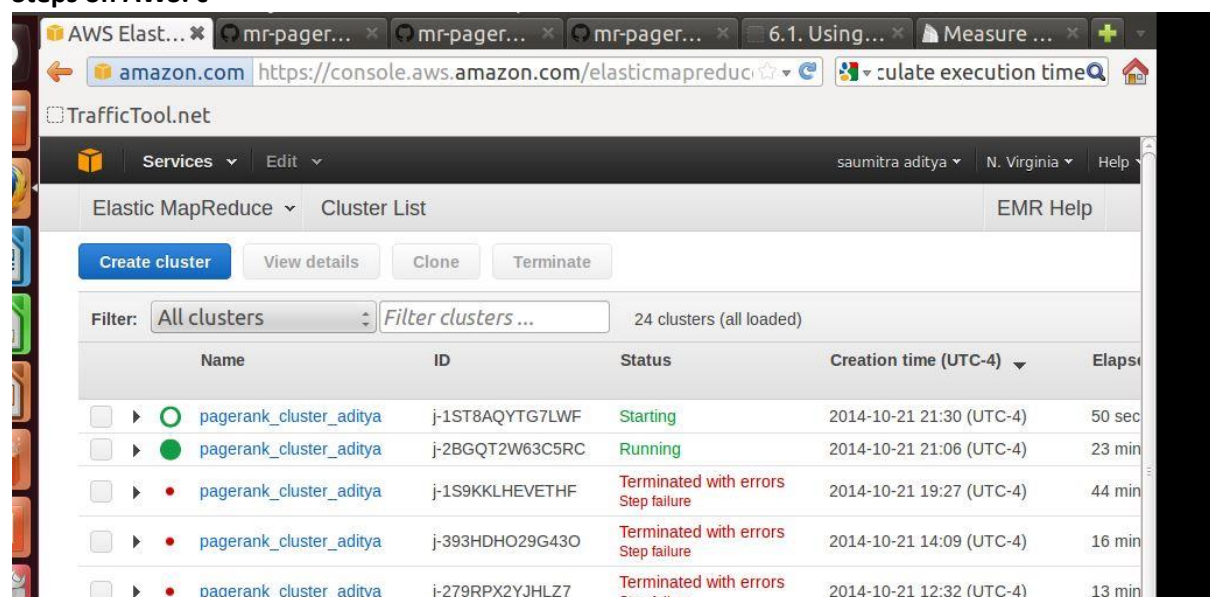
### Method for checking convergence:

I calculate the difference between the updated and last pagerank of a node and sum it over all the nodes, I call this value total\_delta. As one of the arguments I specify a slackness factor eg.01 i.e. I will iterate until sum of delta over all nodes is not less than .01.

### Handling dangling nodes

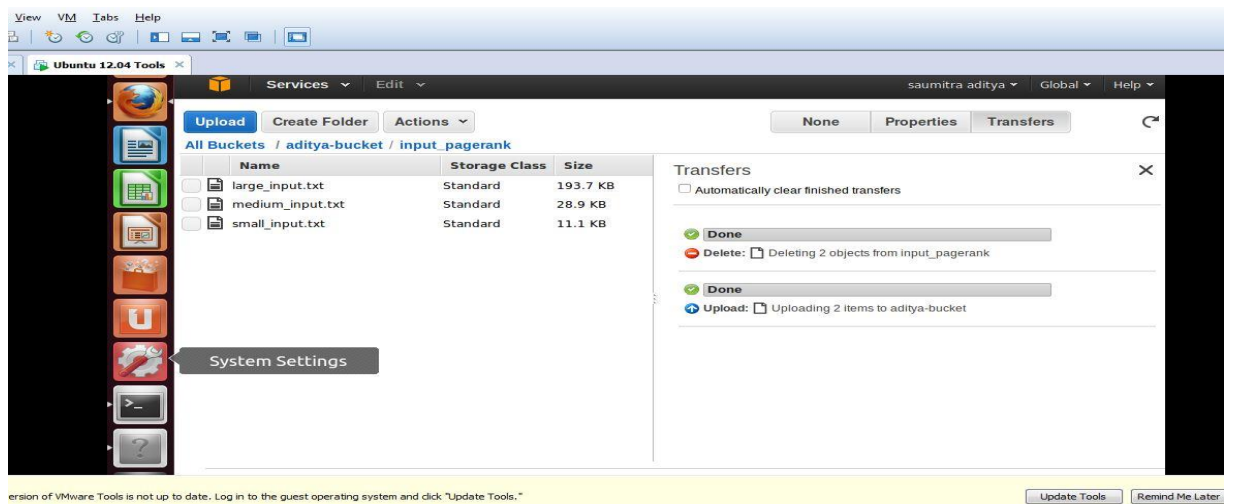
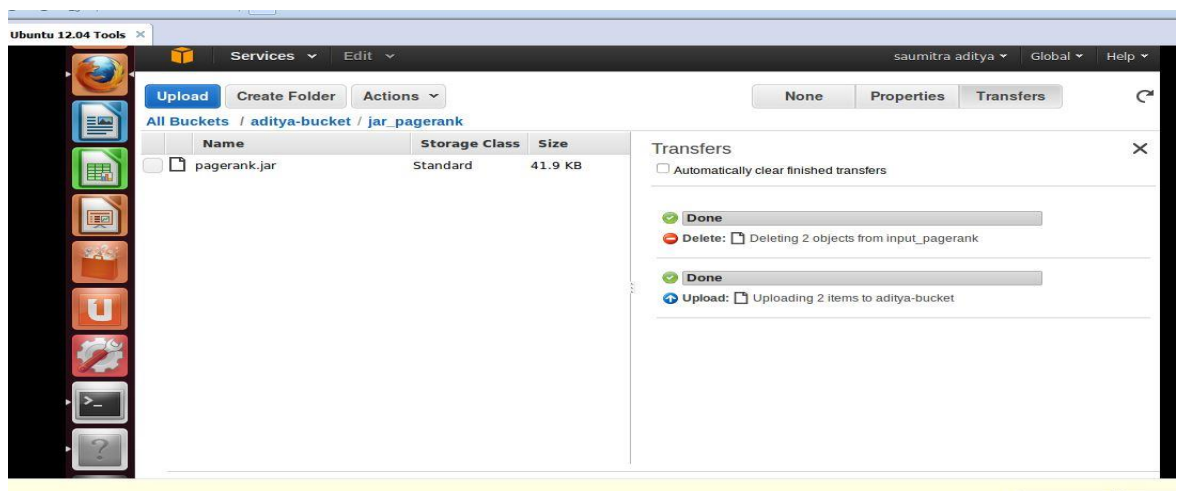
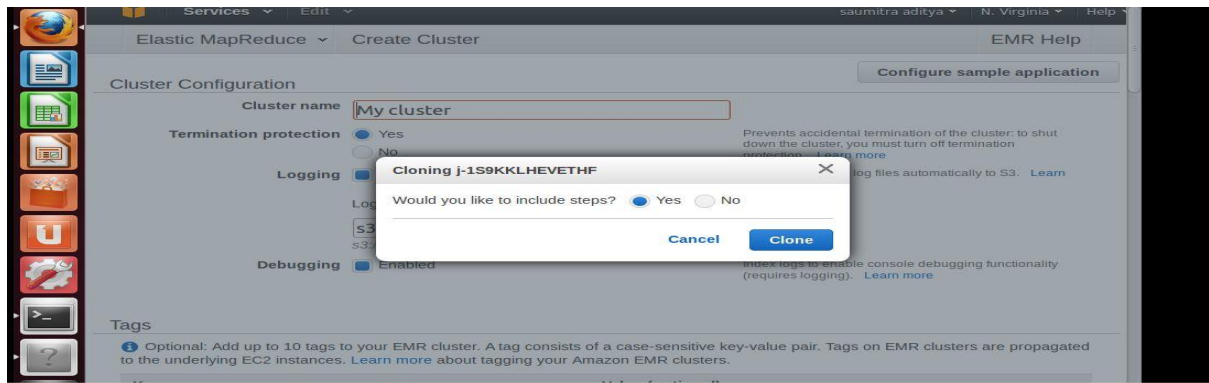
The pagerank of dangling node is distributed among all the nodes in equal proportions i.e every node in the graph gets (pagerank of dangling node/ total nodes) contribution from the dangling node.

### Steps on AWS. c



Name	ID	Status	Creation time (UTC-4)	Elapsed
pagerank_cluster_aditya	j-1ST8AQYT67LWF	Starting	2014-10-21 21:30 (UTC-4)	50 sec
pagerank_cluster_aditya	j-2BGQT2W63C5RC	Running	2014-10-21 21:06 (UTC-4)	23 min
pagerank_cluster_aditya	j-1S9KKLHEVETHF	Terminated with errors Step failure	2014-10-21 19:27 (UTC-4)	44 min
pagerank_cluster_aditya	j-393HDHO29G43O	Terminated with errors Step failure	2014-10-21 14:09 (UTC-4)	16 min
pagerank_cluster_aditya	j-279RPX2YJHLZ7	Terminated with errors Step failure	2014-10-21 12:32 (UTC-4)	13 min

Clone a cluster, upload jar file and input files.



Elastic MapReduce Create Cluster EMR Help

Cluster Configuration [Configure sample application](#)

Cluster name

Termination protection ☐ Yes ☒ No  
Prevents accidental termination of the cluster: to shut down the cluster, you must turn off termination protection. [Learn more](#)

Logging ☒ Enabled  
Copy the cluster's log files automatically to S3. [Learn more](#)

Log folder S3 location  
  
s3://<bucket-name>/<folder>/

Debugging ☒ Enabled  
Index logs to enable console debugging functionality (requires logging). [Learn more](#)

Tags

## Set location of log files.

### Steps

**i** A step is a unit of work you submit to the cluster. A step might contain one or more Hadoop jobs, or contain instructions to install or configure an application. You can submit up to 256 steps to a cluster. [Learn more](#)

Name	Action on failure	JAR location	Arguments
Custom JAR	Terminate cluster	s3://aditya-bucket /jar_pagerank/pagerank.jar	s3://aditya-bucket /input_pagerank /medium_input.txt s3://aditya-bucket/output_pagerank /output_medium 100 .01

Add step

Auto-terminate ☒ Yes ☐ No  
Automatically terminate cluster after the last step is

Specify location of custom jar on S3 and the arguments for the job like input path, output path, maximum iterations and slackness limit.

Ubuntu 12.04 Tools

Services Edit saumitra aditya N. Virginia Help

Elastic MapReduce Cluster List EMR Help

[Create cluster](#) [View details](#) [Clone](#) [Terminate](#)

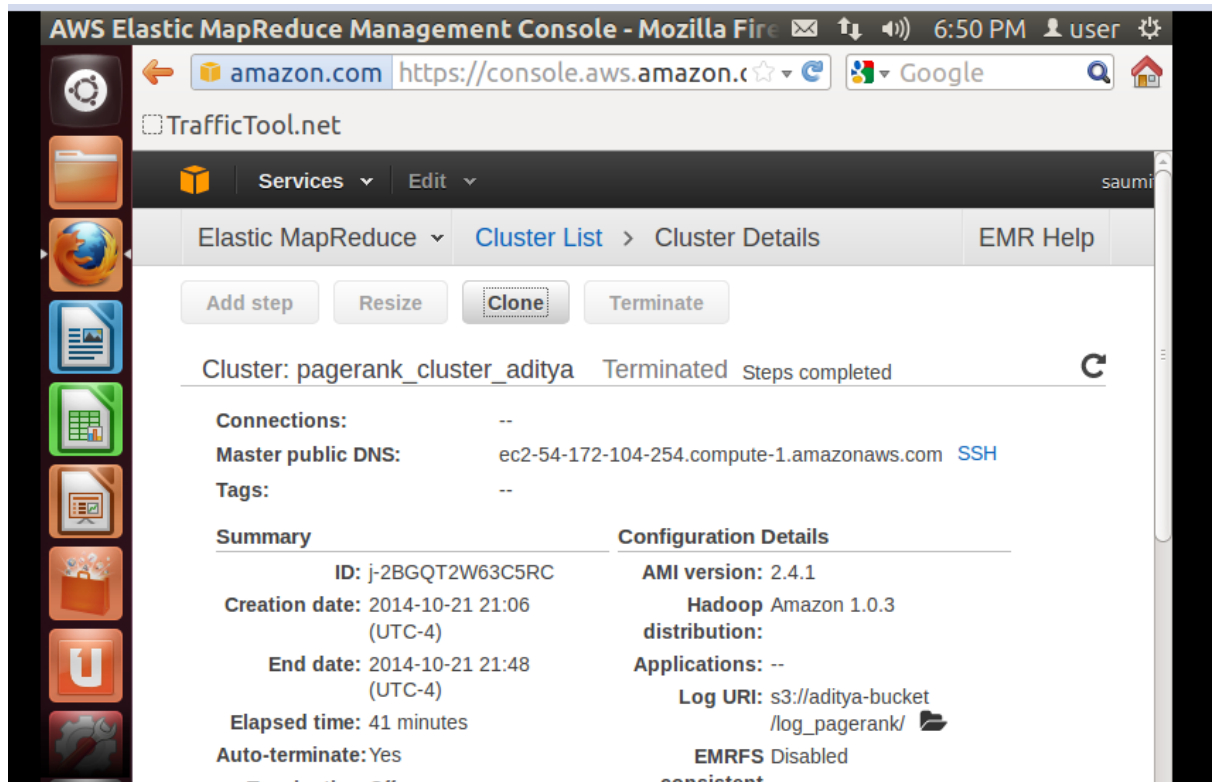
Filter: All clusters Filter clusters ... 23 clusters (all loaded)

Name	ID	Status	Creation time (UTC-4)	Elapsed
<input checked="" type="checkbox"/> pagerank_cluster_aditya	j-2BGQT2W63C5RC	Running	2014-10-21 21:06 (UTC-4)	19 minutes
<input type="checkbox"/> pagerank_cluster_aditya	j-1S9KKLHEVETHE	Terminated with errors Step failure	2014-10-21 19:27 (UTC-4)	44 minutes
<input type="checkbox"/> pagerank_cluster_aditya	j-393HDHO29G43O	Terminated with errors Step failure	2014-10-21 14:09 (UTC-4)	16 minutes
<input type="checkbox"/> pagerank_cluster_aditya	j-279RPX2YJHLZ7	Terminated with errors Step failure	2014-10-21 12:32 (UTC-4)	13 minutes
<input type="checkbox"/> pagerank_cluster_aditya	j-1USF6UX64PRTN	Terminated with errors Step failure	2014-10-21 12:02 (UTC-4)	9 minutes
<input type="checkbox"/> pagerank_cluster_aditya	j-25FD9CM1W43GN	Terminated with errors Step failure	2014-10-20 13:11 (UTC-4)	11 minutes
<input type="checkbox"/> pagerank_cluster_aditya	j-3KOPCIT5N5CBE	Terminated with errors	2014-10-20 13:09 (UTC-4)	25 seconds

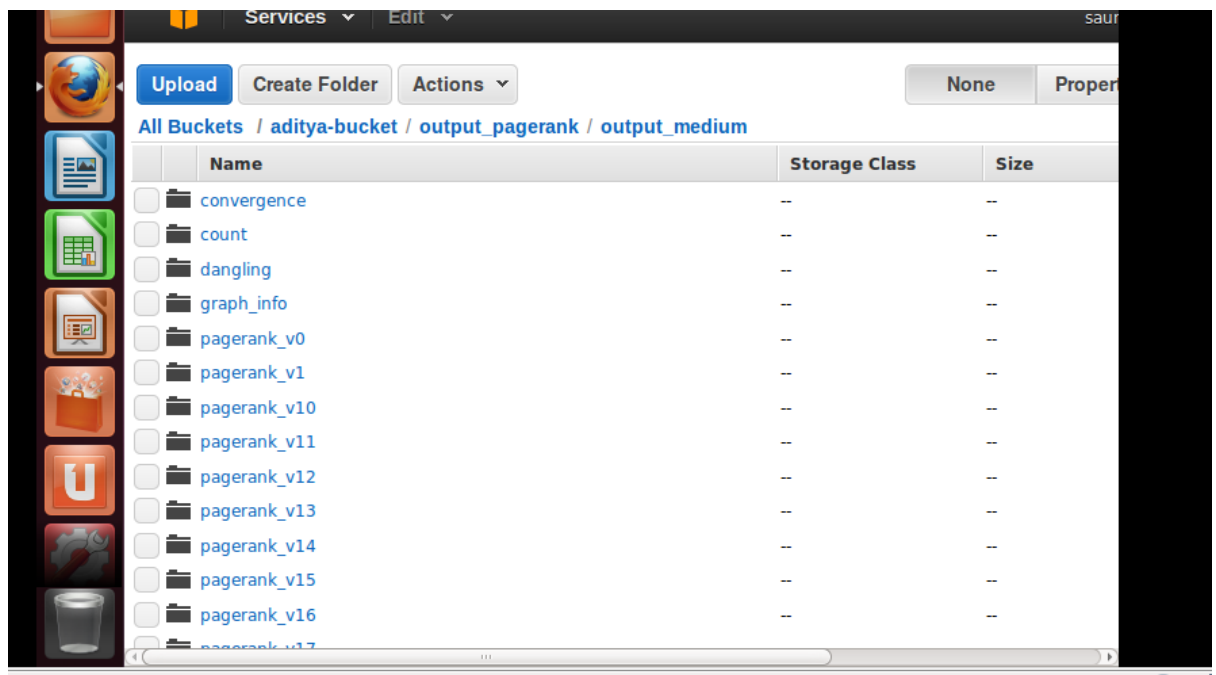
Terminal pagerank\_cluster\_aditya

Latest version of VMware Tools is not up to date. Log in to the guest operating system and click "Update Tools." [Update Tools](#) [Remind Me Later](#)

Start the cluster.



Wait for job to complete.



Output files.

## Results.

### Small input.

graph_summary	total_nodes	93	total_links	195	min_links	0
	max_links	5	avg_links	2.0		

### Top ten nodes

104524212055442757665907965243560045101	0.0032298502535568416
82306156766194587629690350083967473394	0.002391805989823018
134661996234159808488375170070473187582	0.0023640307765620604
155346617108560808581184142629329729230	0.002333236338016026
227109448702302113507903604759918416063	0.002242148643274156
168673252469550579557899067836420932546	0.002120614257327908
57979370615741928609283005034325220088	0.0018038879600958434
17649598783482525745073710167618606107	0.0016266018873613973
97668020538124808838065356267872887871	0.0015550908078008105
284207700310595285051849664560473814989	0.001242114789023274

### Metrics

convergence\_count= 17  
pagerank job completed in 2399015ms

[https://s3.amazonaws.com/aditya-bucket/output\\_pagerank/output\\_small/sorted\\_output/part-r-00000](https://s3.amazonaws.com/aditya-bucket/output_pagerank/output_small/sorted_output/part-r-00000)

[https://s3.amazonaws.com/aditya-bucket/output\\_pagerank/output\\_small/graph\\_info/part-r-00000](https://s3.amazonaws.com/aditya-bucket/output_pagerank/output_small/graph_info/part-r-00000)

### Medium input

graph_summary	total_nodes	316	total_links	430	min_links	0
	max_links	5	avg_links	1.0		

### Top ten nodes

111981443422667599916101641267414970874	0.0013445596494086798
30442676062515284415598723418014355061	7.510357058956637E-4
217182398344717121985059912345853998316	6.580261564678982E-4
64363282148945876210890336872865755343	5.074033228553838E-4
104105844697470013276372331783894076726	4.92268950932983E-4
255141271871887572604204954207769279563	4.900699868634195E-4
116480772629362012002460626777081605400	4.73769838798419E-4
298690743135077500802007851608046438995	4.516332036552296E-4
303806832053566290572095716352649981643	4.390352259543741E-4
Metrv148511838361064104411653673322648403910	4.37235137198252E-4

## Metrics

convergence\_count= 17  
pagerank job completed in 2246659ms

[https://s3.amazonaws.com/aditya-bucket/output\\_pagerank/output\\_medium2/sorted\\_output/part-r-00000](https://s3.amazonaws.com/aditya-bucket/output_pagerank/output_medium2/sorted_output/part-r-00000)

[https://s3.amazonaws.com/aditya-bucket/output\\_pagerank/output\\_medium2/graph\\_info/part-r-00000](https://s3.amazonaws.com/aditya-bucket/output_pagerank/output_medium2/graph_info/part-r-00000)

## Large input.txt

graph_summary	total_nodes	1458	total_links	3545	min_links	0
max_links	5	avg_links	2.0			

## Top ten nodes

64032941963750223601505696787123138445	9.969838542203818E-4
119337412437940133144881923208049882442	8.582750615375713E-4
294418289840301973322672169300394924184	6.918287599413247E-4
1712967822958713490055716528324178036	6.388412220270341E-4
284510239251910046427593057486449185085	5.848073863619843E-4
214917686594559236497547622533457258166	5.628633041528273E-4
131802765080162628666440881374317948788	5.56646529785805E-4
156471617313644419826686265789184402299	5.44236473830064E-4
269011742891273932588164655856253091447	5.358855817756014E-4
259479793941959149960309933682866059975	5.344815197232751E-4

## Metrics

convergence\_count= 17  
pagerank job completed in 2219340ms

[https://s3.amazonaws.com/aditya-bucket/output\\_pagerank/output\\_large2/sorted\\_output/part-r-00000](https://s3.amazonaws.com/aditya-bucket/output_pagerank/output_large2/sorted_output/part-r-00000)

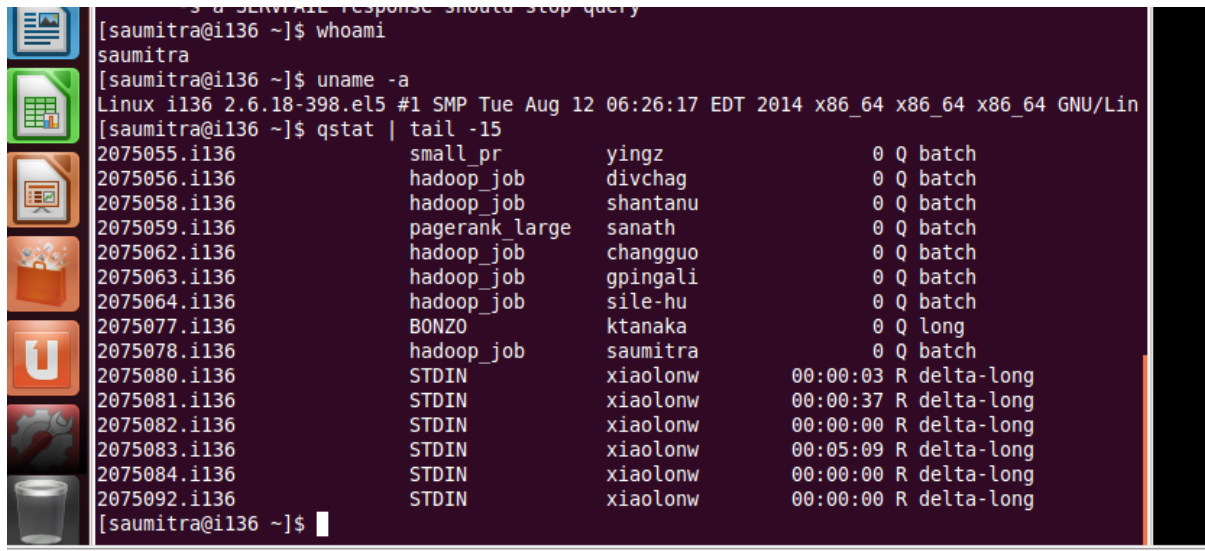
[https://s3.amazonaws.com/aditya-bucket/output\\_pagerank/output\\_large2/graph\\_info/part-r-00000](https://s3.amazonaws.com/aditya-bucket/output_pagerank/output_large2/graph_info/part-r-00000)

## Conclusion:

The metrics indicate that all the three graphs converge after 17 iterations when slackness factor is .01 (it is the sum of delta of latest and last PR over all nodes.)

The timing metric from AWS is not accurate, while computing on local machine larger the input size more is the running time.

I was not able to run jobs on Future Grid as it is still in queue after waiting for more than a day.

A terminal window with a dark purple background and white text. On the left side, there is a vertical dock with several application icons: a blue document icon, a green spreadsheet icon, an orange folder icon, a red folder icon, an orange 'U' icon, a red and white wrench icon, and a glass icon. The terminal text shows the user 'saumitra' at host 'i136' performing several commands. The 'uname -a' command displays system details including the kernel version '2.6.18-398.el5' and architecture 'x86\_64'. The 'qstat | tail -15' command shows a list of 15 jobs. The first 8 jobs are in a 'batch' state, and the last 7 jobs are in a 'delta-long' state.

```
[saumitra@i136 ~]$ whoami
saumitra
[saumitra@i136 ~]$ uname -a
Linux i136 2.6.18-398.el5 #1 SMP Tue Aug 12 06:26:17 EDT 2014 x86_64 x86_64 x86_64 GNU/Linux
[saumitra@i136 ~]$ qstat | tail -15
2075055.i136      small_pr      yingz         0 Q batch
2075056.i136      hadoop_job    divchag       0 Q batch
2075058.i136      hadoop_job    shantanu      0 Q batch
2075059.i136      pagerank_large sanath         0 Q batch
2075062.i136      hadoop_job    changguo      0 Q batch
2075063.i136      hadoop_job    gpingali      0 Q batch
2075064.i136      hadoop_job    sile-hu       0 Q batch
2075077.i136      BONZO        ktanaka       0 Q long
2075078.i136      hadoop_job    saumitra      0 Q batch
2075080.i136      STDIN        xiaolonw      00:00:03 R delta-long
2075081.i136      STDIN        xiaolonw      00:00:37 R delta-long
2075082.i136      STDIN        xiaolonw      00:00:00 R delta-long
2075083.i136      STDIN        xiaolonw      00:05:09 R delta-long
2075084.i136      STDIN        xiaolonw      00:00:00 R delta-long
2075092.i136      STDIN        xiaolonw      00:00:00 R delta-long
[saumitra@i136 ~]$
```